

Efficient MinHash-based Algorithms for Big Structured Data



Wei Wu

Faculty of Engineering and Information Technology
University of Technology Sydney

A thesis is submitted for the degree of
Doctor of Philosophy

July 2018

I would like to dedicate this thesis to my loving parents ...

Declaration

I hereby declare that except where specific reference is made to the work of others, the contents of this dissertation are original and have not been submitted in whole or in part for consideration for any other degree or qualification in this, or any other University. This dissertation is the result of my own work and includes nothing which is the outcome of work done in collaboration, except where specifically indicated in the text. This dissertation contains less than 65,000 words including appendices, bibliography, footnotes, tables and equations and has less than 150 figures.

Wei Wu
July 2018

Acknowledgements

I would like to express my earnest thanks to my principal supervisor, Dr. Ling Chen, and co-supervisor, Dr. Bin Li, who have provided the tremendous support and guidance for my research. Dr. Ling Chen always controlled the general direction of my PhD life. She made the recommendation letter for me, and contacted the world leading researchers for collaboration with me. Also, she encouraged me to attend the research conferences, where my horizon was remarkably broaden and I knew more about the research frontiers by discussing with the worldwide researchers. Dr. Bin Li guided me to explore the unknown, how to think and how to write. He once said to me, "It is not that I am supervising you. Instead, this is a spiraling process and we make progress together." In the first two years of PhD, I was still groping in my research topic and did not obtain any results. During that time, I was very anxious and unconfident. However, my supervisors were patient enough and always encouraged me. No doubt that I cannot acquire my current research performance without their patient supervision.

I am grateful to my best friend, Dr. Chuan Luo, who shares three-year research life in Peking University. Although we are currently in different research communities, we often discuss the research trends. In my most depressing days, he enlightened me patiently, even though we were in different countries.

I am grateful to all members in CAI for their participation in and valuable comments on my research. Particularly, I thank Prof. Chengqi Zhang, Dr. Chunyang Liu, Dr. Sujuan Hou, Dr. Meng Fang, Dr. Qinzhe Zhang, Dr. Weiwei Liu, Mr. Yu Liu, Mr. Bo Han, Mr. Yuangang Pan, Mr. Jiangchao Yao, Mr. Adi Lin, Mr. Daokun Zhang, Mr. Dong Wen, Dr. Qian Zhang and Miss Jiamiao Wang, and other students for their help in my PhD study.

I am grateful to the financial support from UTS scholarship and FEIT Travel funds.

Finally, and above all, I would like to thank my parents for their unconditional support and encouragement from emotion and finance. They always believe in me, give me invaluable suggestions and fully support all my decisions. No words could possibly express my deepest gratitude for their endless love, self-sacrifice and unwavering help.

To them I dedicate the dissertation.

Abstract

Nowadays, data are growing explosively in the information society. Every day, an incredibly large number of data are generated from social networks, financial industries, medical device and digital equipment, etc. – Big data have been driving data mining research in both academia and industry. No matter how data mining and machine learning develop, in most tasks such as classification, clustering and retrieval, it is one of the most fundamental operations to compute the similarity of data. However, exact similarity computation has become daunting for big data due to the “3V” nature (volume, velocity and variety). Therefore, the thesis aims to develop a number of efficient randomized hashing algorithms for big structured data to approximate the similarity.

Taking into account the popularity and importance of data modelled as weighted sets, graphs/networks, the MinHash-based algorithms are proposed to transform weighted sets, graphs and network nodes into low-dimensional hash codes for efficient similarity estimation.

In Chapter 4, we propose two MinHash-based algorithms for weighted sets, which are two instances of the Consistent Weighted Sampling (CWS) scheme, for real-value weighted sets by uncovering the working mechanism of the state-of-the-art weighted MinHash algorithm, Improved Consistent Weighted Sampling (ICWS) algorithm. The two algorithms both represent a weighted set as a low-dimensional hash code. The first algorithm, the Canonical Consistent Weighted Sampling (CCWS) algorithm, reconstructs the hash functions in ICWS in order to overcome the flaw that ICWS breaks the uniformity of the MinHash scheme. Consequently, our CCWS algorithm not only shares the same theoretical computational complexity as ICWS but also strictly complies with the definition of the MinHash algorithm. The second algorithm, the Practical Consistent Weighted Sampling (PCWS) algorithm, simplifies ICWS by transforming the original form of ICWS into an equivalent expression. As a result, our PCWS algorithm is not only mathematically equivalent to ICWS and preserves the same theoretical properties, but also saves 20% memory footprint and substantial computational cost compared to ICWS.

In Chapter 5, we propose a MinHash-based algorithm for graphs, i.e., the K -Ary Tree Hashing (KATH) algorithm, for representing a graph as a low-dimensional feature vector.

The main idea of KATH is to construct a traversal table to quickly approximate the subtree patterns in Weisfeiler-Lehman (WL) graph kernel using K -ary trees. Based on the traversal table, our KATH algorithm employs a recursive indexing process that performs only r times of matrix indexing to generate all $(r - 1)$ -depth K -ary trees, where the leaf node labels of a tree can uniquely specify the pattern. After that, the MinHash scheme is used to fingerprint the acquired subtree patterns for a graph. The experimental results show that our KATH algorithm runs significantly faster than state-of-the-art methods while achieving competitive or better accuracy.

In Chapter 6, we propose a MinHash-based algorithm for network nodes, i.e., the NetHash algorithm. Consequently, each node in the network is mapped to a low-dimensional vector space. Our NetHash algorithm employs the randomized hashing technique to encode shallow trees, each of which is rooted at a node of the network. The main idea is to efficiently encode both attributes and structure information of each node by recursively sketching the corresponding rooted tree from bottom (i.e., highest-order neighboring nodes) to top (i.e., root node), and particularly, to preserve as much information closer to the root node as possible. The extensive experimental results show that the proposed algorithm, which does not need learning, runs significantly faster than the state-of-the-art learning-based network embedding methods while achieving competitive or even better performance in accuracy. Furthermore, the simplicity of the algorithm enables our NetHash algorithm to embed a millions-of-nodes network within 1,000 seconds on a single machine.

Table of contents

List of figures	xv
List of tables	xix
Nomenclature	xix
1 Introduction	1
1.1 Motivation	1
1.2 Problem Statement	2
1.2.1 The MinHash-based Algorithms for Weighted Sets	2
1.2.2 The MinHash-based Algorithm for Graphs	3
1.2.3 The MinHash-based Algorithm for Network Nodes	4
1.3 Contributions	5
1.4 Structure	7
1.5 Publications	8
2 Preliminary	9
2.1 Definitions	9
2.2 Notations	10
2.3 Benchmark Data Sets	11
2.3.1 Weighted Sets	11
2.3.2 Graphs	13
2.3.3 Networks	16
3 Literature Review	17
3.1 Hashing Techniques	17
3.2 Similarity Hashing	18
3.2.1 Learning to Hash	20
3.2.2 Locality Sensitive Hashing	22

3.3	MinHash	23
3.4	Structured Hashing	25
4	Two MinHash-based Algorithms for Weighted Sets	27
4.1	Introduction	27
4.2	Preliminaries	31
4.2.1	MinHash	31
4.2.2	Weighted MinHash	32
4.2.3	Integer Weighted MinHash	33
4.3	Review of Improved CWS	34
4.4	Canonical Consistent Weighted Sampling	37
4.4.1	Breach of Uniformity	37
4.4.2	The CCWS Algorithm	39
4.4.3	Proof of Correctness	41
4.4.4	Remarks	43
4.5	Practical Consistent Weighted Sampling	44
4.5.1	ICWS Revisited	44
4.5.2	The PCWS Algorithm	46
4.5.3	Analysis	47
4.6	Experimental Results	50
4.6.1	Experimental Preliminaries	50
4.6.2	Experiments for CCWS	51
4.6.3	Experiments for PCWS	55
4.7	Related Work	61
4.8	Conclusion	63
5	A MinHash-based Algorithm for Graphs	65
5.1	Introduction	65
5.2	Preliminaries	67
5.2.1	Problem Description	68
5.2.2	Weisfeiler-Lehman Graph Kernels	68
5.2.3	Nested Subtree Hash Kernels	70
5.2.4	MinHash	71
5.3	K -ary Tree Hashing	72
5.3.1	The Algorithm	74
5.3.2	Complexity Analysis	79
5.4	Experiments	79

5.4.1	Experimental Preliminaries	80
5.4.2	Data Sets	81
5.4.3	Mix	82
5.4.4	Results on NCI	84
5.4.5	Results on qHTS	85
5.4.6	Results on NCI-Yeast	86
5.4.7	Results on DBLP	87
5.4.8	Results on Twitter	88
5.4.9	Results on 10K Synthetic Graph Sets	89
5.4.10	Impact of Traversal Table Size on KATH	89
5.5	Related Work	90
5.6	Conclusion	92
6	A MinHash-based Algorithm for Network Nodes	93
6.1	Introduction	93
6.2	Preliminaries	96
6.2.1	Problem Definition	96
6.2.2	MinHash	96
6.3	NetHash	97
6.3.1	The NetHash Algorithm	98
6.3.2	Exponentially Decayed Diffusion	99
6.3.3	Theoretical Analysis	101
6.4	Experimental Results	103
6.4.1	Parameter Settings	105
6.4.2	Node Classification on Citation Networks	106
6.4.3	Link Prediction on Social Networks	108
6.4.4	Case Study: Node Retrieval on a Large-scale Citation Network . . .	108
6.4.5	Discussion on the Results	109
6.5	Related Work	110
6.5.1	Network Node Embedding	110
6.5.2	Graph Hashing	111
6.6	Conclusion	112
7	Conclusions and Future Work	113
7.1	Summary of This Thesis	113
7.2	Future Work	114

References

115

List of figures

1.1	Comparison between a binary set and a weighted set for a document.	3
1.2	Comparison between a graph and a set for a chemical compound, CH_3Cl	4
1.3	A citation network.	5
1.4	Thesis structure.	7
3.1	A toy example for a traditional hash function and collision.	18
3.2	Comparison between a traditional hash function and the similarity hashing scheme.	19
3.3	Categorization for Similarity Hashing.	21
3.4	An illustration for MinHash.	24
4.1	An illustration of “active indices”. The red dashed lines represent “active indices” whose hash values are monotonically decreasing from bottom to top; while the blue dashed lines are non-active subelements, where $v_{k,y_4} < v_{k,y_1} < \{v_{k,y_2}, v_{k,y_3}\}$. IWMH proceeds traversing “active indices” from bottom to top by skipping many non-active subelements.	33
4.2	Breach of uniformity. The left bars represent the original weights, while the right ones take the logarithm of the original weights. The red arrows show that the samples are mapped from the original weights onto the logarithmic ones. Because of the logarithmic transform, the samples in different subelements of \mathcal{T} and \mathcal{R} are mapped into the same subelement, which increases the collision probability of the “active indices”.	38
4.3	Performance results, classification accuracy (left column) and runtime (right column) of the compared methods, on the four real-world text data sets. The x -axis denotes the length of fingerprints, D	53
4.4	Performance results, classification accuracy (left column) and runtime (right column) of the compared methods, on the four real-world text data sets. The x -axis denotes the length of fingerprints, D	54

4.5	Classification results in accuracy (left column) and runtime (right column) of the compared methods on Real-sim, Rcv1 and KDD. The x -axis denotes the length of fingerprints, D	56
4.6	Retrieval results in Precision@ K (first columns) and MAP@ K (second columns) of the compared methods on Webspam. The x -axis denotes the length of fingerprints, D	58
4.7	Retrieval results in Precision@ K (first columns) and MAP@ K (second columns) of the compared methods on Webspam. The x -axis denotes the length of fingerprints, D	59
4.8	Retrieval results in Precision@ K (first columns) and MAP@ K (second columns) of the compared methods on Url. The x -axis denotes the length of fingerprints, D	60
4.9	Retrieval results in Precision@ K (first columns) and MAP@ K (second columns) of the compared methods on Url. The x -axis denotes the length of fingerprints, D	61
4.10	Retrieval runtime of the compared methods on Webspam and Url. The x -axis denotes the length of fingerprints, D	62
5.1	Illustration of the WL kernel and the NSH kernel with $r = 2$	69
5.2	Two examples of traversal table construction. The first column of the traversal tables are the five root nodes which are followed by their neighboring nodes sorted or min-hashed in terms of labels.	72
5.3	Classification Accuracy (upper row) and CPU Time (bottom row) comparison results on the three graph data sets in the group of Real-World Mix. The x -axis denotes the depth of subtrees.	83
5.4	Classification Accuracy (odd rows) and CPU Time (even rows) comparison results on the nine graph data sets in NCI. The x -axis denotes the depth of subtrees.	84
5.5	Classification Accuracy (odd rows) and CPU Time (even rows) comparison results on the nine graph data sets in NCI. The x -axis denotes the depth of subtrees.	85
5.6	Classification Accuracy (odd rows) and CPU Time (even rows) comparison results on the nine graph data sets in NCI. The x -axis denotes the depth of subtrees.	86
5.7	Classification Accuracy (left) and CPU Time (right) comparison results on qHTS. The x -axis denotes the depth of subtrees.	87

5.8	Classification Accuracy (left) and CPU Time (right) comparison results on NCI-Yeast. The x -axis denotes the depth of subtrees.	87
5.9	Classification Accuracy (left) and CPU Time (right) comparison results on DBLP. The x -axis denotes the depth of subtrees.	88
5.10	Classification Accuracy (left) and CPU Time (right) comparison results on Twitter. The x -axis denotes the depth of subtrees.	88
5.11	Classification Accuracy and CPU Time comparison results on Synthetic 10K20L10E (KN short for KATH-naive, and KM short for KATH-MinHash).	89
5.12	Classification Accuracy comparison results on Real-World NCI1. The x -axis denotes the size of traversal table. The left figure represents KATH-naive, and the right figure represents KATH-MinHash.	90
6.1	An illustration of embedding a node (recursively sketching a rooted tree) on a network.	98

List of tables

2.1	Important notations used in the thesis	10
2.2	Summary of the six weighted set data sets	12
2.3	Summary of the seventeen graph data sets.	14
2.4	Summary of the seven network data sets.	16
5.1	Summary of distribution of No. neighbors for NCI.	90
6.1	Tuning depth for Cora.	104
6.2	Tuning decay rate for Cora.	104
6.3	Node classification results on citation networks.	105
6.4	Link prediction results on social networks.	106
6.5	Top-5 Query results on a large-scale citation network.	107

